

# Autoencoders and Autoregressive Models

## Deep Learning (DSE316/616)

Vinod K Kurmi  
*Assistant Professor, DSE*

Indian Institute of Science Education and Research Bhopal

Oct 13, 2022

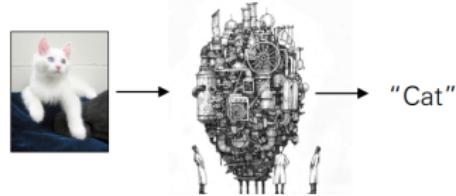


# Disclaimer

- Much of the material and slides for this lecture were borrowed from
  - Bernhard Schölkopf's MLSS 2017 lecture,
  - Tommi Jaakkola's 6.867 class,
  - CMP784: Deep Learning Fall 2021 Erkut Erdem Hacettepe University
  - Fei-Fei Li, Andrej Karpathy and Justin Johnson's CS231n class
  - Hongsheng Li's ELEG5491 class
  - Tsz-Chiu Au slides
  - Mitesh Khapra Class notes

# Discriminative vs. Generative Models

$$p(y|x)$$



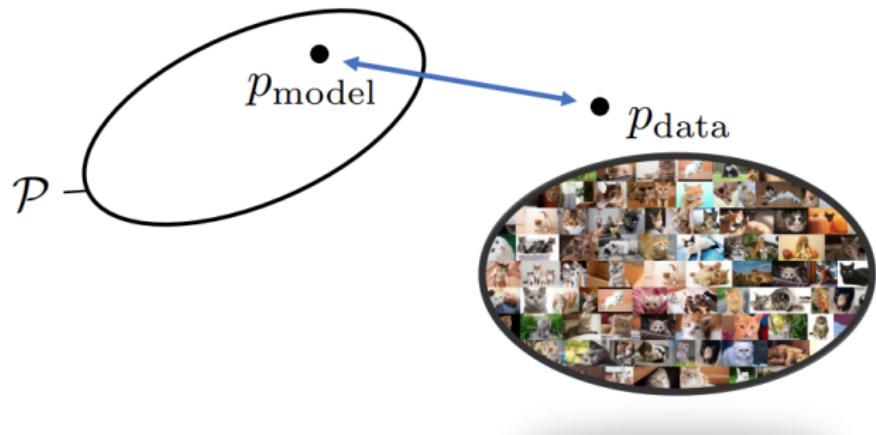
Discriminative models

$$p(x|y)$$



Generative models

# Generative Modeling



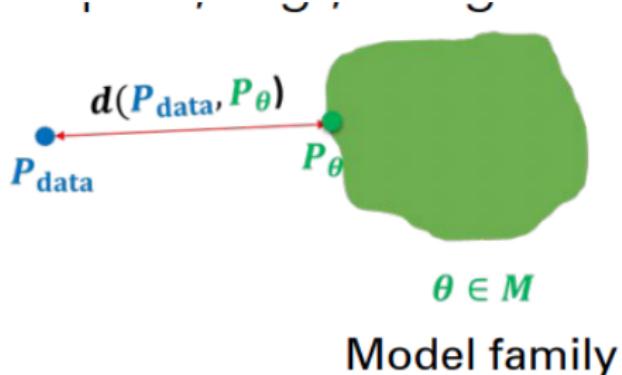
- **Goal:** Learn some underlying hidden structure of the training samples to generate novel samples from same data distribution

# Learning a generative model

- We are given a training set of examples, e.g., images of dogs

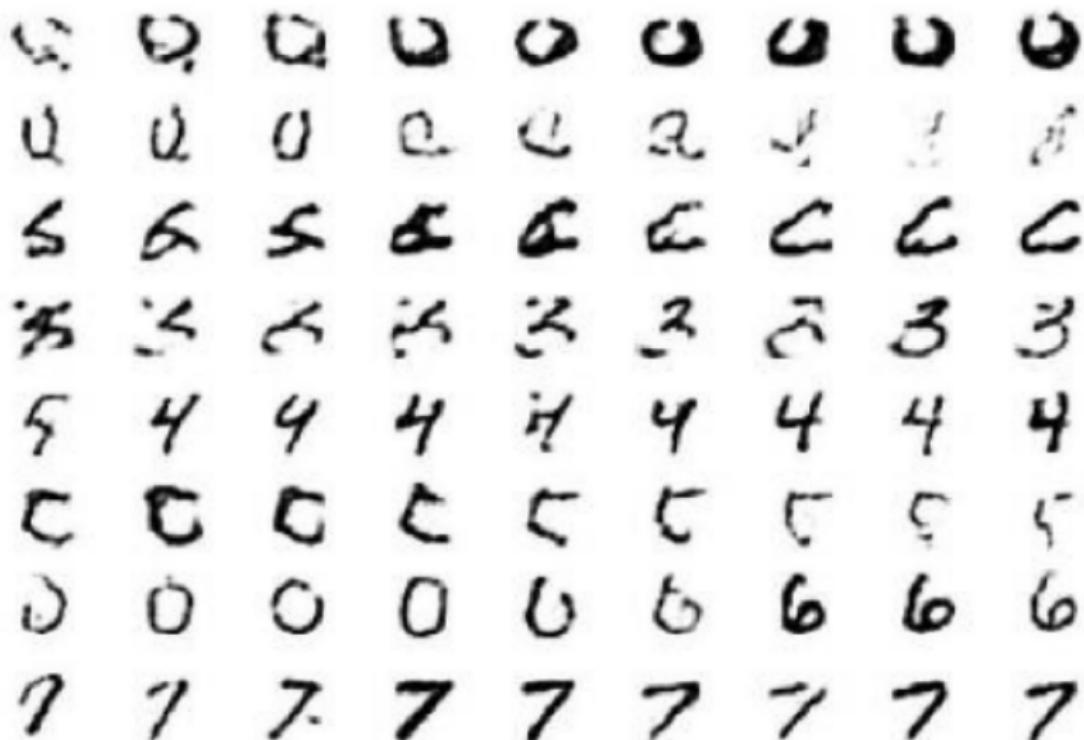


$$\mathbf{x}_i \sim P_{\text{data}} \\ i = 1, 2, \dots, n$$



- We want to learn a probability distribution  $p(x)$  over images  $x$  s.t.
  - Generation:** If we sample  $x_{\text{new}} \sim p(x)$ ,  $x_{\text{new}}$  should look like a dog (sampling)
  - Density estimation:**  $p(x)$  should be high if  $x$  looks like a dog, and low otherwise (anomaly detection)
  - Unsupervised representation learning:** We should be able to learn what these images have in common, e.g., ears, tail, etc. (features)

## Generate Images



# Generate Images

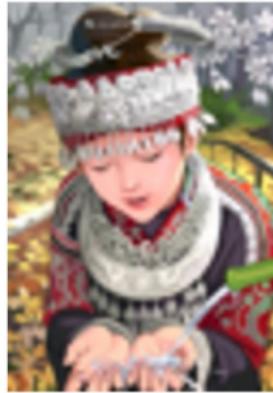


# Generate Images



# Generate Images

bicubic  
(21.59dB/0.6423)



SRResNet  
(23.53dB/0.7832)



SRGAN  
(21.15dB/0.6868)



original



# Generate Images



# Generate Images



# Generate Audio



1 Second



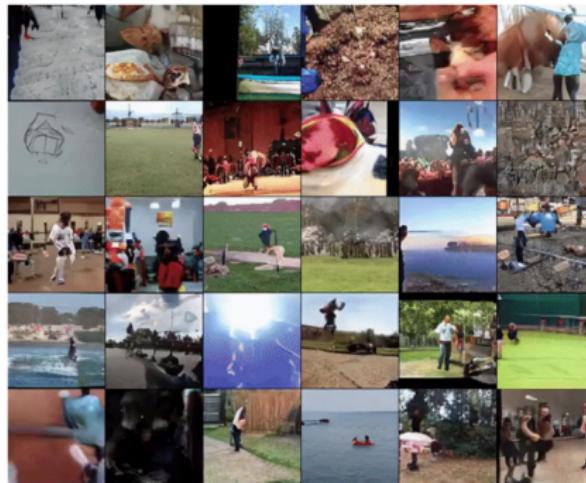
Parametric



WaveNet



# Generate Video



DVD-GAN: Adversarial Video Generation on Complex Datasets, Clark, Donahue, Simonyan, 2019

# Generate Text

PANDARUS:

Alas, I think he shall be come approached and the day  
When little strain would be attain'd into being never fed,  
And who is but a chain and subjects of his death,  
I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,  
Breaking and strongly should be buried, when I perish  
The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

[Char-rnn, karpathy, 2015]

# Generate Math

```
\begin{proof}
We may assume that $\mathcal{I}$ is an abelian sheaf on
$\mathcal{C}$.
\item Given a morphism $\Delta : \mathcal{F} \rightarrow \mathcal{I}$ is an injective and let $\mathfrak{q}$ be an abelian sheaf on $X$.
Let $\mathcal{F}$ be a fibered complex. Let $\mathcal{F}$ be a category.
\begin{enumerate}
\item \hyperref[setain-construction-phantom]{Lemma}
\label{lemma-characterize-quasi-finite}
Let $\mathcal{F}$ be an abelian quasi-coherent sheaf on
$\mathcal{C}$.
Let $\mathcal{F}$ be a coherent $\mathcal{O}_X$-module. Then
$\mathcal{F}$ is an abelian catenary over $\mathcal{C}$.
\item The following are equivalent
\begin{enumerate}
\item $\mathcal{F}$ is an $\mathcal{O}_X$-module.
\end{enumerate}
\end{enumerate}
\end{proof}
```

For  $\bigoplus_{i=1,\dots,n} U_i = 0$ , where  $\mathcal{L}_{i,i} = 0$ , hence we can find a closed subset  $H$  in  $H$  and any sets  $\mathcal{F}$  on  $X$ ,  $U$  is a closed immersion of  $S$ , then  $U \rightarrow T$  is a separated algebraic space.

*Proof.* Proof of (1). It also start we get

$$S = \text{Spec}(R) = U \times_X U \times_X U$$

and the compatibility in the fibre product covering we have to prove the lemma generated by  $\coprod Z \times_U U \rightarrow V$ . Consider the maps  $M$  along the set of points  $\text{Sch}_{/\text{fppf}}$  and  $U \rightarrow U'$  is the fibre category of  $S$  in  $U$  in Section ?? and the fact that any  $U$  affine, see Morphisms, Lemma ?? Hence we obtain a scheme  $S$  and any open subset  $W \subset U$  in  $\text{Sh}(G)$  such that  $\text{Spec}(R') \rightarrow S$  is smooth or an

$$U = \bigcup U_i \times_S U_i$$

which has a nonzero morphism we may assume that  $f_i$  is of finite presentation over  $S$ . We claim that  $\mathcal{O}_{X,i}$  is a scheme where  $x, x', x'' \in S'$  such that  $\mathcal{O}_{X,x'} \rightarrow \mathcal{O}_{X,x''}$  is separated. By Algebra, Lemma ?? we can define a map of complexes  $\text{GL}_{S'}(x'/S'')$  and we win.  $\square$

To prove study we see that  $\mathcal{F}_{|U_i}$  is a covering of  $X^i$ , and  $T_i$  is an object of  $\mathcal{F}_{X^i/S}$  for  $i > 0$  and  $\mathcal{F}_0$  exists and let  $\mathcal{F}_i$  be a presheaf of  $\mathcal{O}_X$ -modules on  $C$  as a  $\mathcal{F}$ -module. In particular  $\mathcal{F} = U/\mathcal{F}$  we have to show that

$$\tilde{M}^\bullet = \mathbb{Z}^\bullet \otimes_{\text{Spec}(k)} \mathcal{O}_{X,i} - i_X^{-1} \mathcal{F}$$

is a unique morphism of algebraic stacks. Note that

$$\text{Arrows} = (\text{Sch}/S)_{/\text{fppf}}^{\text{op}}, (\text{Sch}/S)_{/\text{fppf}}$$

and

$$V = \Gamma(S, \mathcal{O}) \longrightarrow (U, \text{Spec}(A))$$

is an open subset of  $X$ . Thus  $U$  is affine. This is a continuous map of  $X$  is the inverse, the groupoid scheme  $S$ .

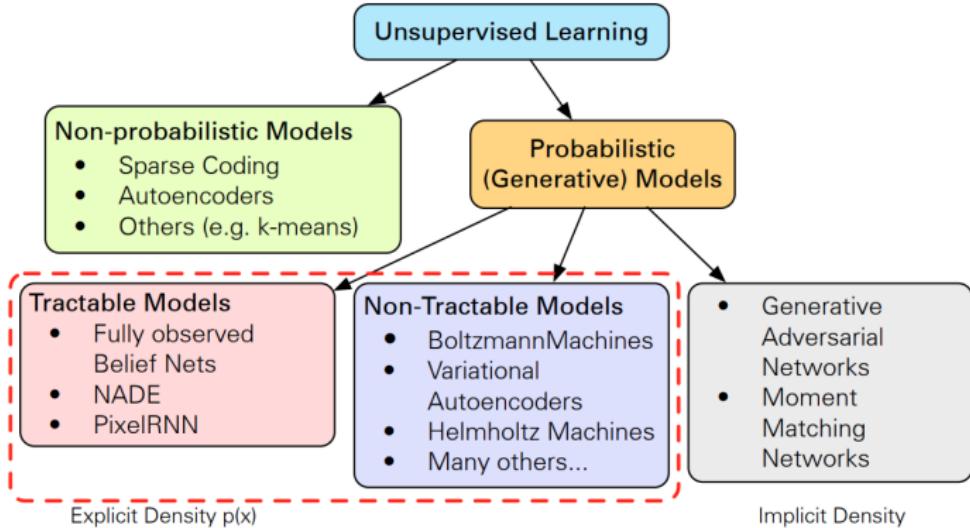
*Proof.* See discussion of sheaves of sets.  $\square$

The result for prove any open covering follows from the less of Example ?? It may replace  $S$  by  $X_{\text{space}, \text{state}}$  which gives an open subspace of  $X$  and  $T$  equal to  $S_{2, \sigma}$ , see Descent, Lemma ?? Namely, by Lemma ?? we see that  $R$  is geometrically regular over  $S$ .

[Char-rnn, karpathy, 2015]

# Why Unsupervised Learning?

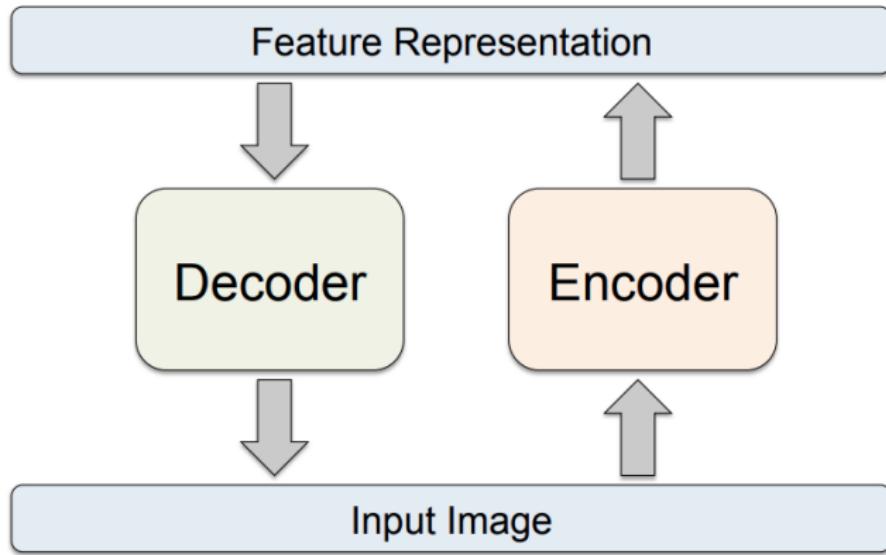
- Given high-dimensional data , we want to find a low-dimensional model characterizing the population.
- Recent progress mostly in supervised DL
- Real challenges for unsupervised DL
- Potential benefits:
- Exploit tons of unlabeled data
  - Answer new questions about the variables observed
  - Regularizer – transfer learning – domain adaptation
  - Easier optimization (divide and conquer)
  - Joint (structured) outputs



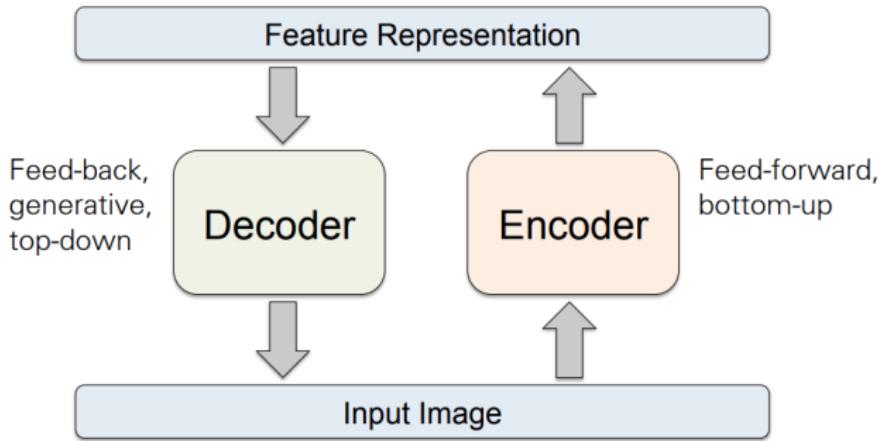
## Basic Building Blocks:

- Autoencoders
- Variational Autoencoders
- Generative Adversarial Networks
- Normalizing Flow Models
- Diffusion models

# Autoencoder

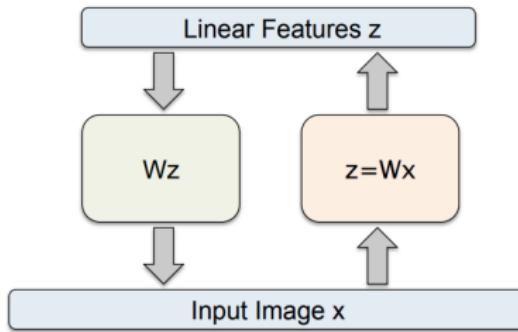


# Autoencoder



- Details of what goes inside the encoder and decoder matter!
- Need constraints to avoid learning an identity

# Autoencoder



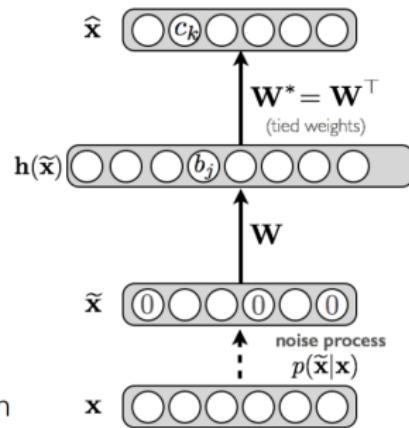
- With nonlinear hidden units, we have a nonlinear generalization of PCA.
- If the **hidden and output layers are linear**, it will learn hidden units that are a linear function of the data and minimize the squared error.
- The K hidden units will span the same space as the first k principal components. The weight vectors may not be orthogonal.

# Denoising Autoencoder

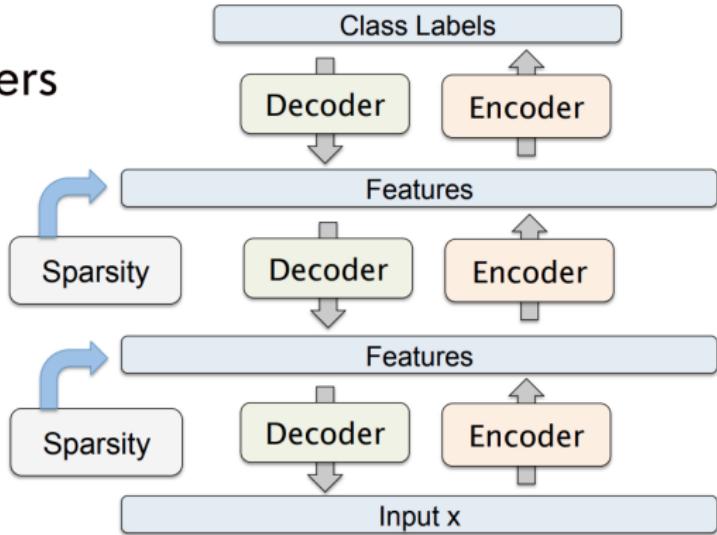
- **Idea:** Representation should be robust to introduction of noise:

- random assignment of subset of inputs to 0, with probability  $\nu$
- Similar to dropouts on the input layer
- Gaussian additive noise

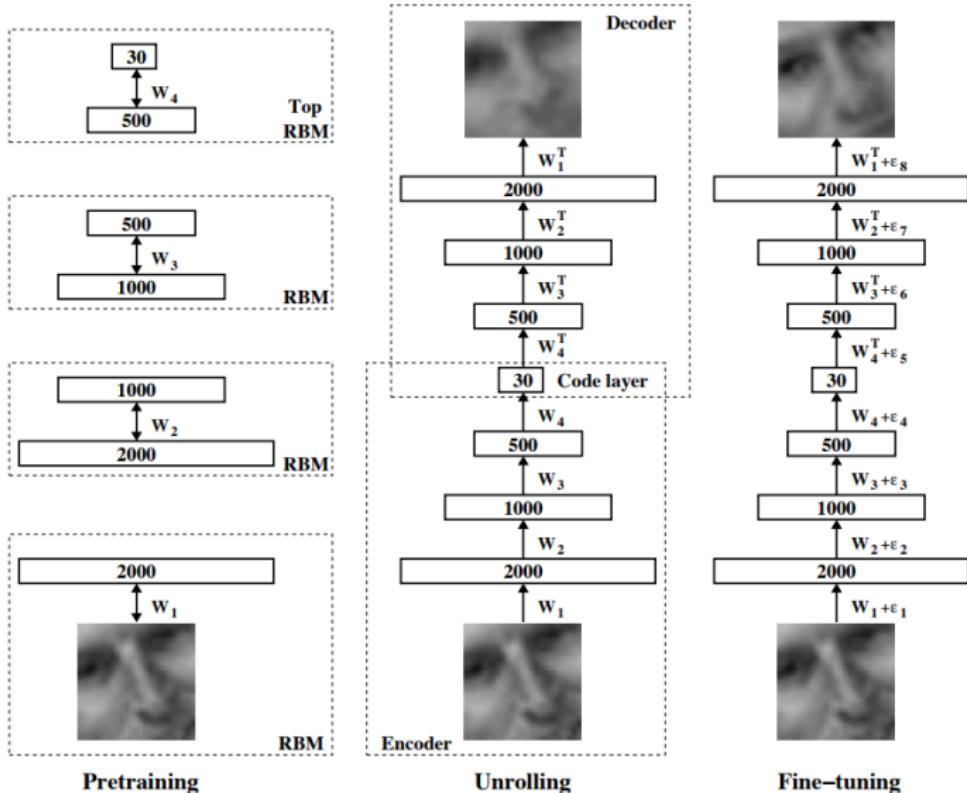
- **Reconstruction**  $\hat{\mathbf{x}}$  computed from the corrupted input  $\tilde{\mathbf{x}}$
- **Loss function** compares  $\hat{\mathbf{x}}$  reconstruction with the noiseless input  $\mathbf{x}$



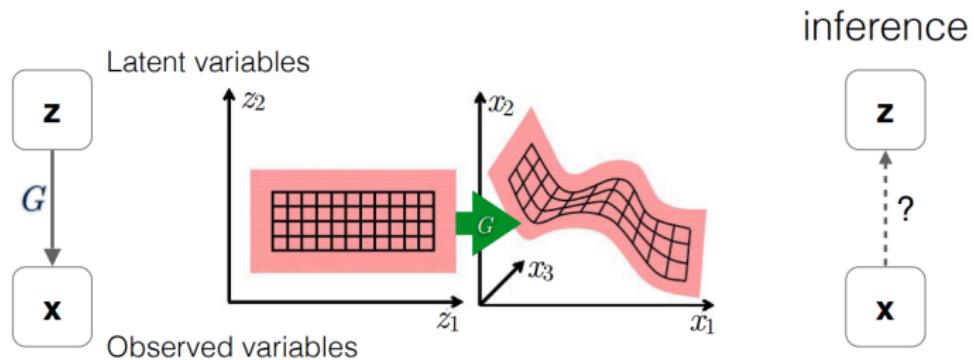
# Stacked Autoencoders



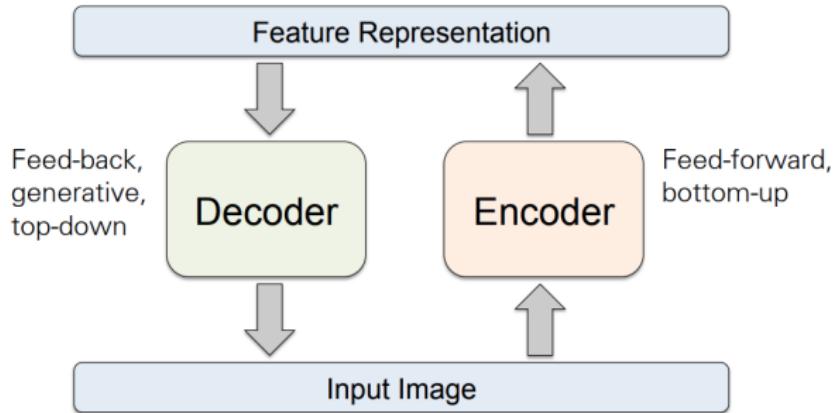
# Deep Autoencoders



# Latent variable model



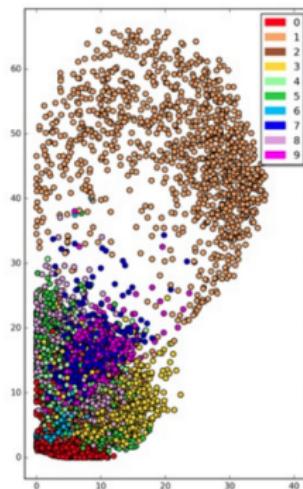
# Autoencoders



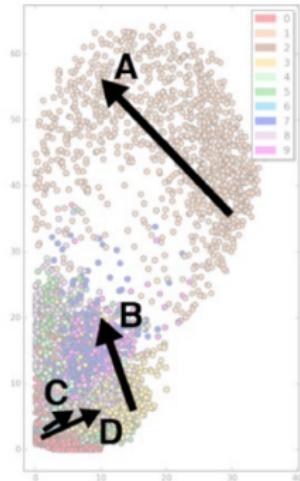
- Details of what goes inside the encoder and decoder matter!
- Need constraints to avoid learning an identity.

# Parameter space of autoencoder

- Let's examine the latent space of an AE.
- Is there any separation of the different classes? If the AE learned the "essence" of the MNIST images, similar images should be close to each other.
- Plot the latent space and examine the separation.
- Here we plot the 2 PCA components of the latent space.



# Traversing the latent space

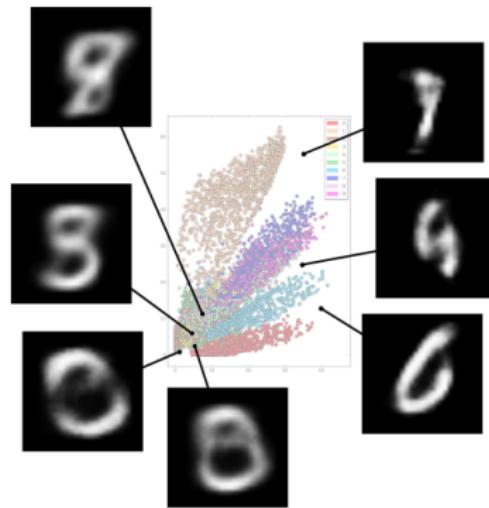


- We start at the start of the arrows in latent space and then move to end of the arrow in 7 steps.
- For each value of z we use the already trained decoder to produce an image.



# Problems with Autoencoders

- Gaps in the latent space
- Discrete latent space
- Separability in the latent space



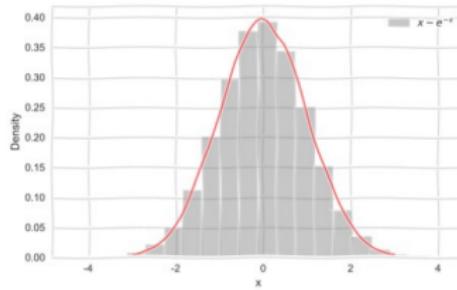
# Generative models

- Imagine we want to generate data from a distribution,

$$x \sim p(x)$$

- e.g.

$$x \sim N(\mu, \sigma)$$



# Generative models

- In other words we can think that if we choose  $z \sim Uniform$  then there is a mapping:

$$x = f(z)$$

such as

$$x \sim p(x)$$

- where in general  $f$  is some complicated function.

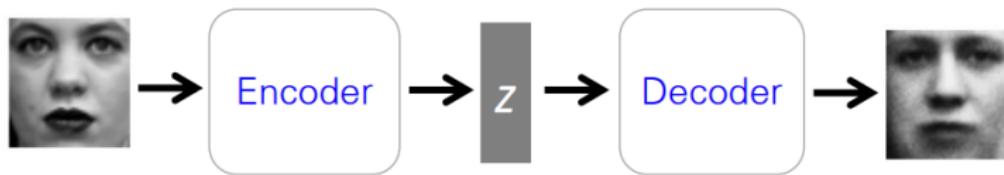
$$x \sim N(\mu, \sigma)$$

- We already know that Neural Networks are great in learning complex functions

$$z \sim g(z) \longrightarrow x = f(z) \longrightarrow x \sim p(x)$$

# Traditional Autoencoders

- In traditional autoencoders, we can think of encoder and decoders as some function mapping

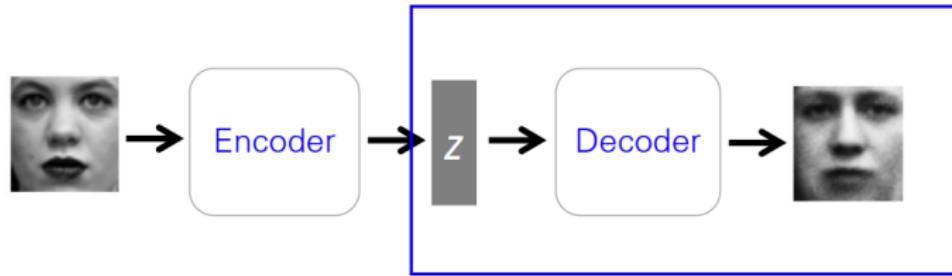


$$z = h(x)$$

$$\hat{x} = f(z)$$

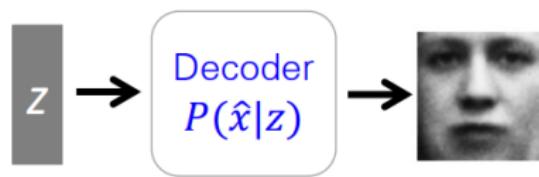
# Variational Autoencoders

- To go to variational autoencoders, we need to first add some stochasticity and think of it as a probabilistic modeling.



# Variational Autoencoders

Sample from  $g(z)$   
e.g. Standard Gaussian

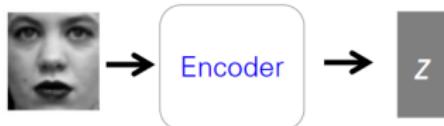


$$z \sim g(z)$$

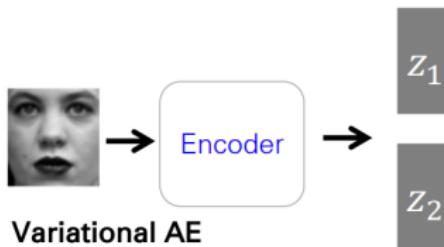
$$\hat{x} = f(z)$$

$$\hat{x} \sim P(x|z)$$

# Variational Autoencoders



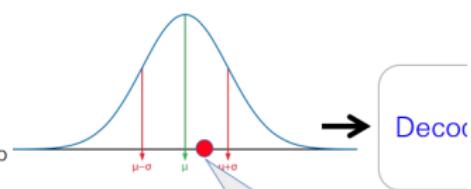
Traditional AE



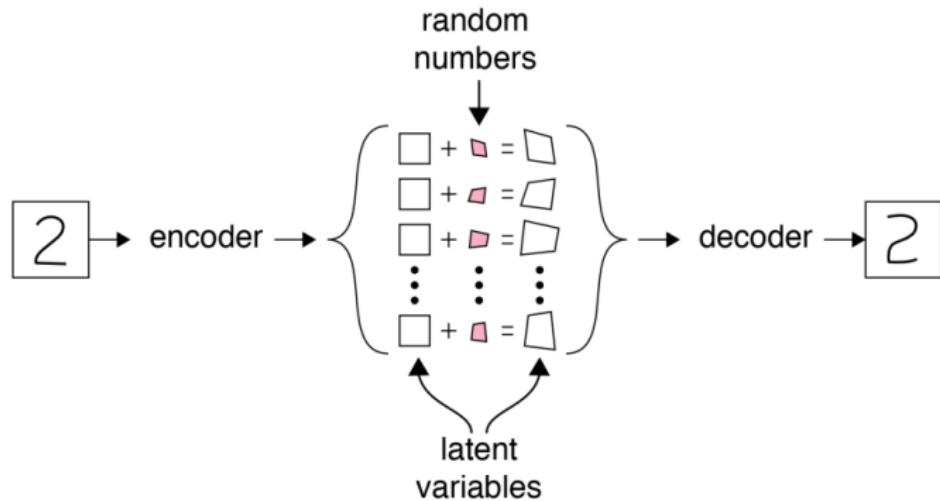
Variational AE

Consider this  
to be the mean  
of a normal  $\mu$

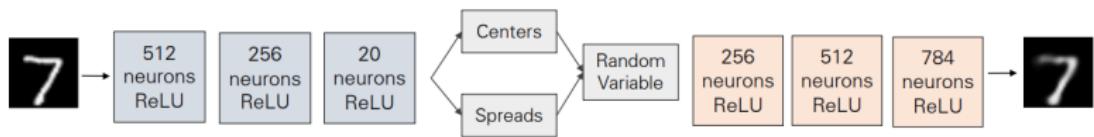
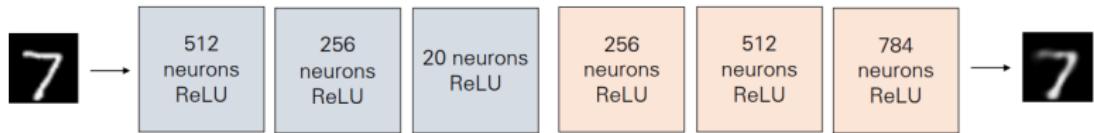
Consider this to  
be the std of a  
normal  $\sigma$



# Variational Autoencoders

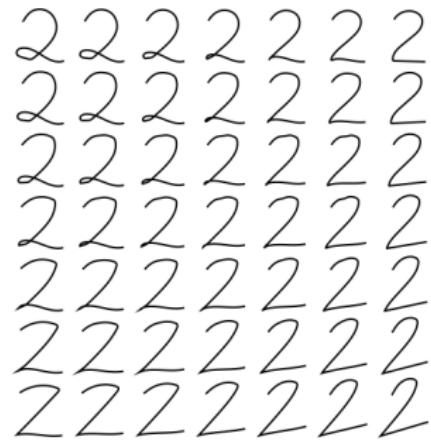


# Variational Autoencoders

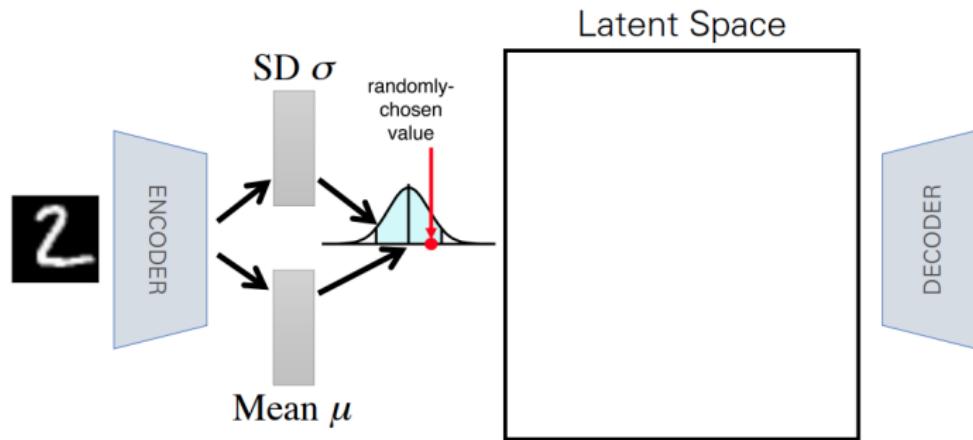


# Separability in Variational Autoencoders

- Separability is not only between classes but we also want similar items in the same class to be near each other.
- For example, there are different ways of writing “2”, we want similar styles to end up near each other.
- Let’s examine VAE, there is something magic happening once we add stochasticity in the latent space.

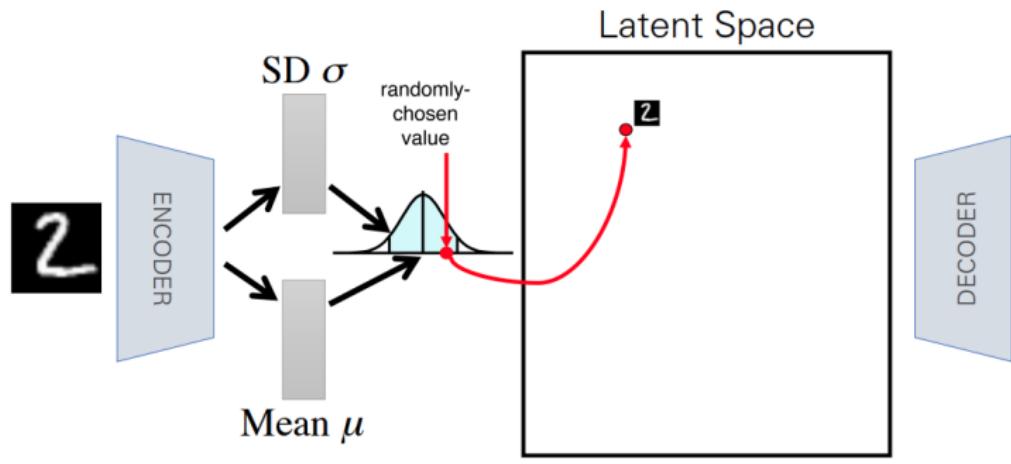


# Separability in Variational Autoencoders



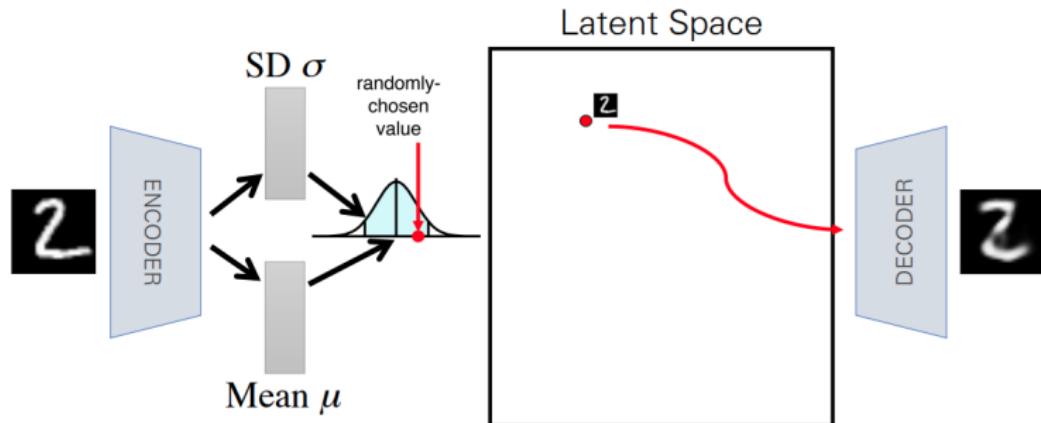
Encode the first sample (a "2") and find  $\mu_1, \sigma_1$

# Separability in Variational Autoencoders



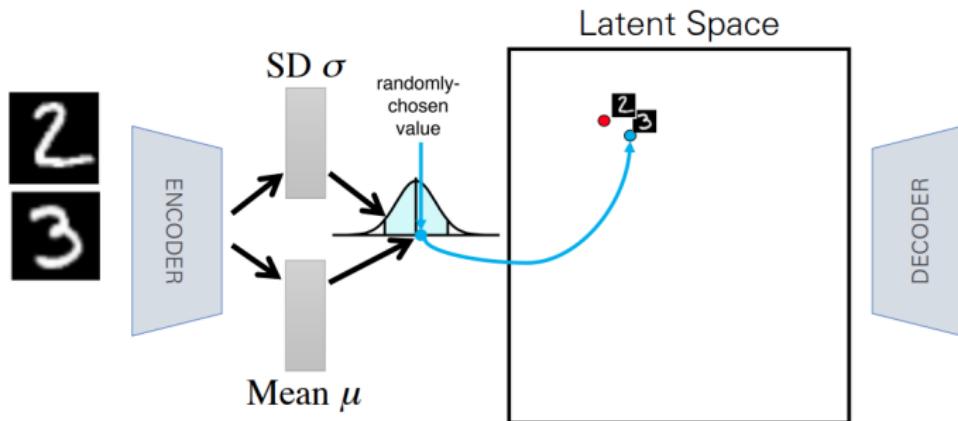
Sample  $z_1 \sim N(\mu_1, \sigma_1)$

# Separability in Variational Autoencoders



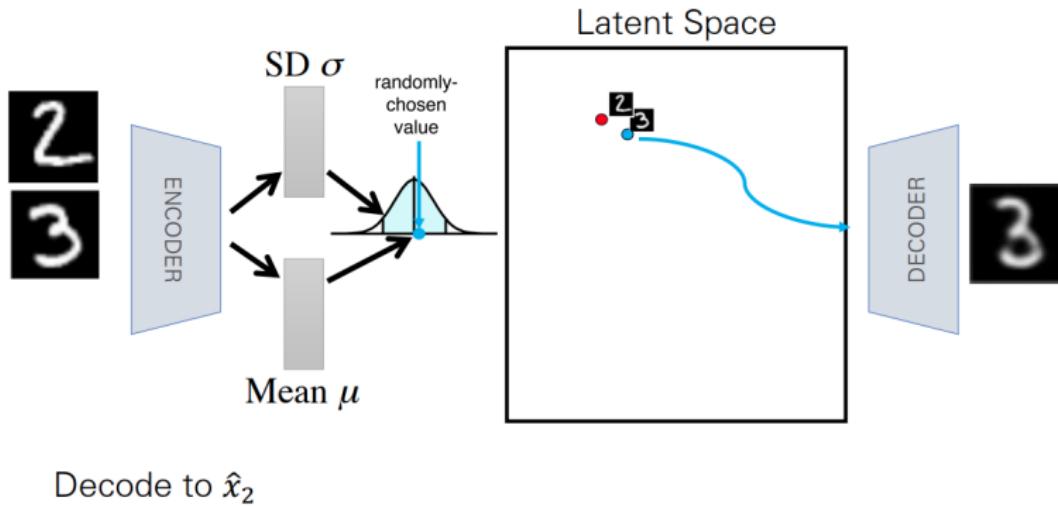
Decode to  $\hat{x}_1$

# Separability in Variational Autoencoders

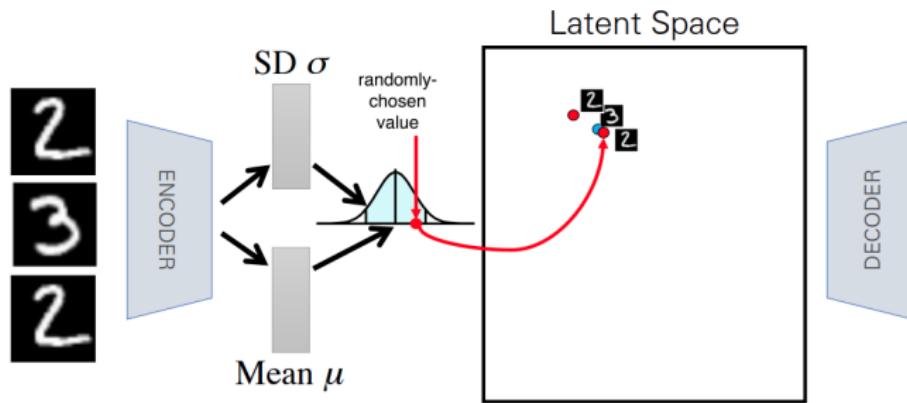


Encode the second sample (a "3") find  $\mu_2, \sigma_2$ . Sample  $z_2 \sim N(\mu_2, \sigma_2)$

# Separability in Variational Autoencoders

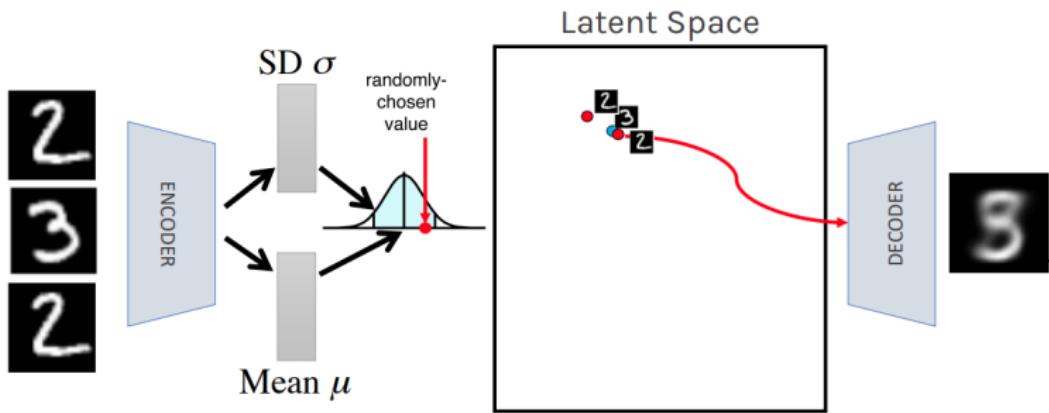


# Separability in Variational Autoencoders



Train with the first sample (a "2") again and find  $\mu_1, \sigma_1$ . However  $z_1 \sim N(\mu_1, \sigma_1)$  will not be the same. It can happen to be close to the "3" in latent space.

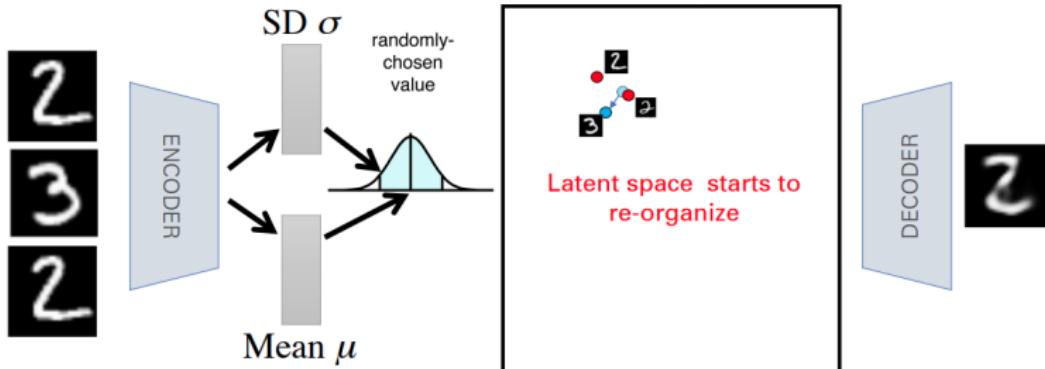
# Separability in Variational Autoencoders



Decode to  $\hat{x}_1$ . Since the decoder only knows how to map from latent space to  $\hat{x}$  space, it will return a "3".

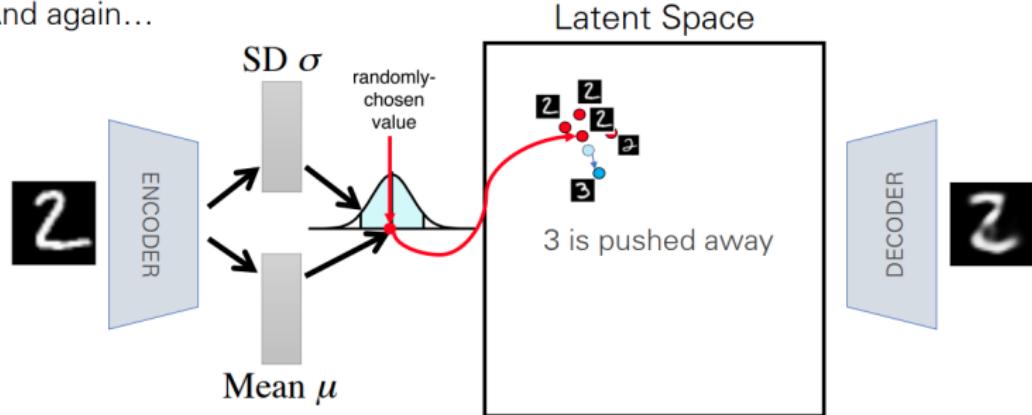
# Separability in Variational Autoencoders

Train with 1<sup>st</sup> sample again



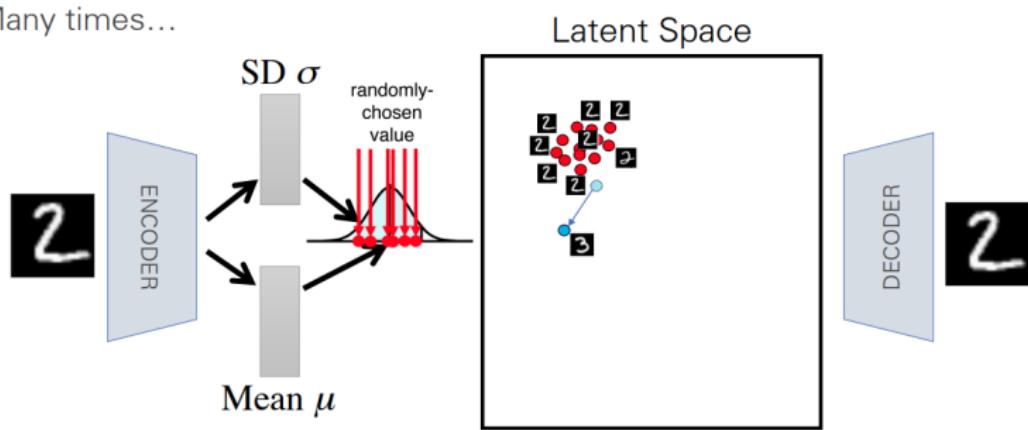
# Separability in Variational Autoencoders

And again...



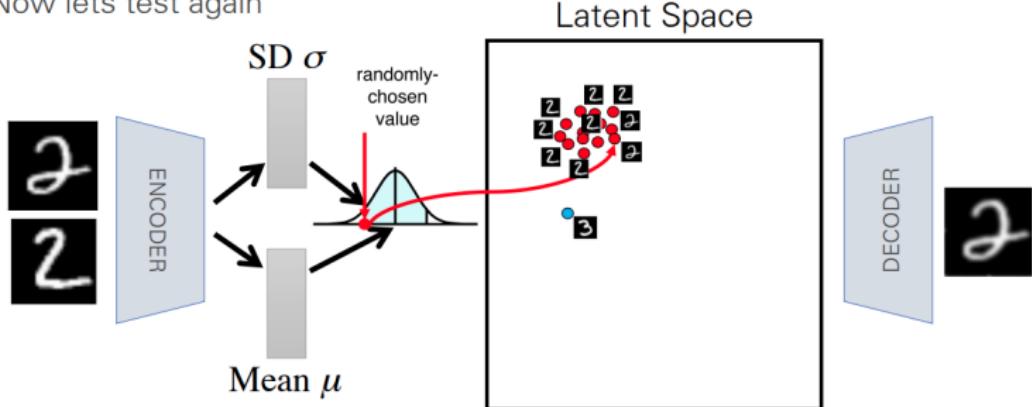
# Separability in Variational Autoencoders

Many times...



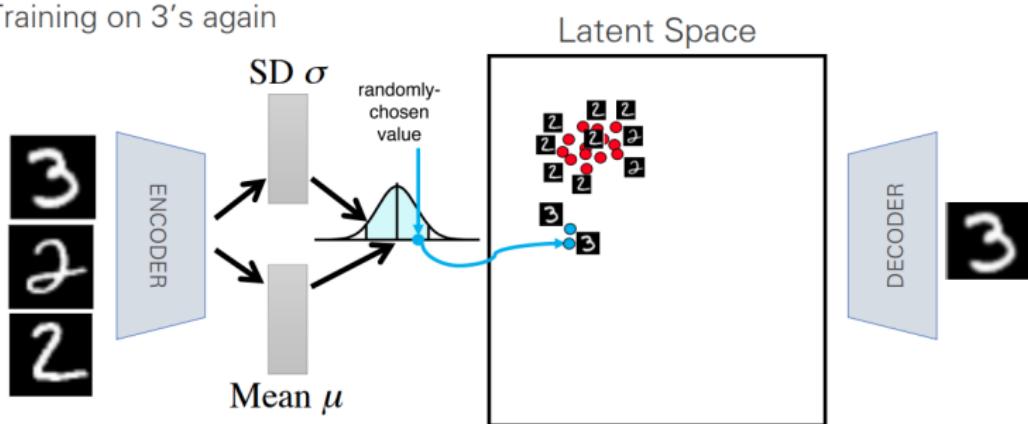
# Separability in Variational Autoencoders

Now lets test again



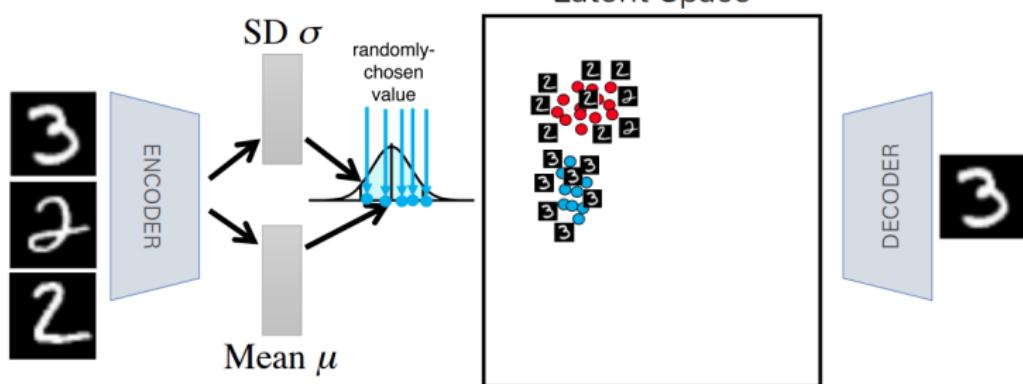
# Separability in Variational Autoencoders

Training on 3's again

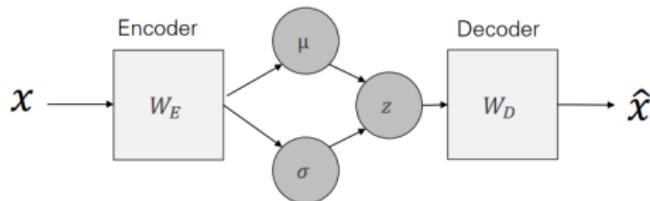


# Separability in Variational Autoencoders

Many times...



# Separability in Variational Autoencoders



Training means learning  $W_E$  and  $W_D$ .

- Define a loss function  $\mathcal{L}$
- Use stochastic gradient descent (or Adam) to minimize  $\mathcal{L}$

The Loss function:

- Reconstruction error:  $\mathcal{L}_R = \frac{1}{n} \sum_i (x_i - \hat{x}_i)^2$
- Similarity between the probability of  $z$  given  $x$ ,  $p(z|x)$ , and some predefined probability distribution  $p(z)$ , which can be computed by Kullback-Leibler divergence (KL):  
$$KL(p(z|x)||p(z))$$

# Training Variational Autoencoders

**Problem:**  $z$  is the dimensionality of your latent space, which can be too large. In other words this  $\int p(\hat{x}|z, x)p(z|x)dz$  becomes intractable.

Instead we turn this into a minimization problem – Variational Calculus  
Find a  $q(z|x)$  that is similar to  $p(z|x)$  by minimizing their difference.

After some math:

$$\text{Reconstruction Loss} \quad \begin{array}{l} \text{Proposal distribution} \\ \text{should resemble} \\ \text{a Gaussian} \end{array}$$
$$-\mathbf{E}_{z \sim q_\phi(z|x)} \log(p_\theta(x|z)) + KL(q_\phi(z|x) \| p_\theta(z))$$

Evidence Lower  
BOund (ELBO)

# Variational AE

- The VAE approach: introduce an inference machine  $q_\phi(z | x)$  that learns to approximate the posterior  $p_\theta(z | x)$ .
  - Define a variational lower bound on the data likelihood:  $p_\theta(x) \geq \mathcal{L}(\theta, \phi, x)$

$$\begin{aligned}\mathcal{L}(\theta, \phi, x) &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z | x)] \\ &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z) + \log p_\theta(z) - \log q_\phi(z | x)] \\ &= -D_{\text{KL}}(q_\phi(z | x) \| p_\theta(z)) + \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)]\end{aligned}$$

*regularization term*      *reconstruction term*

- What is  $q_\phi(z | x)$  ?

# Variational AE

$$\begin{aligned}\log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[ \log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \underbrace{\mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi) \text{ "Elbow"}} + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))}_{\geq 0} \end{aligned}$$

$$\log p_\theta(x^{(i)}) \geq \mathcal{L}(x^{(i)}, \theta, \phi)$$

Variational lower bound (elbow)

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

Training: Maximize lower bound

55

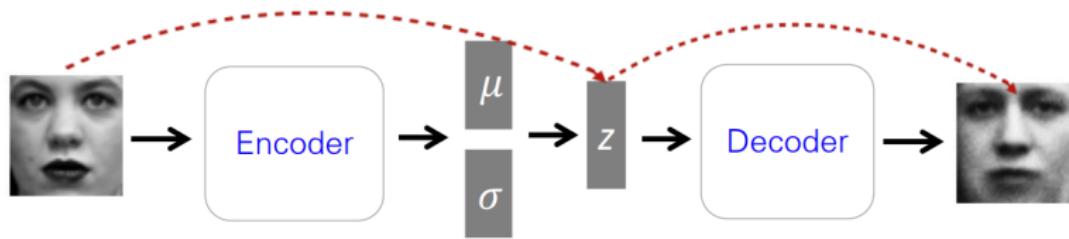
# Training VAE

- Apply stochastic gradient descent (SGD)

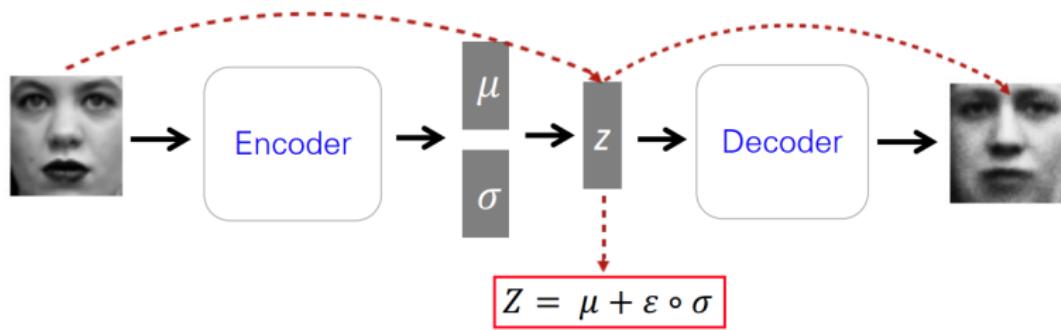
Problem:

- Sampling step not differentiable
- Use a re-parameterization trick
  - Move sampling to input layer, so that the sampling step is independent of the model

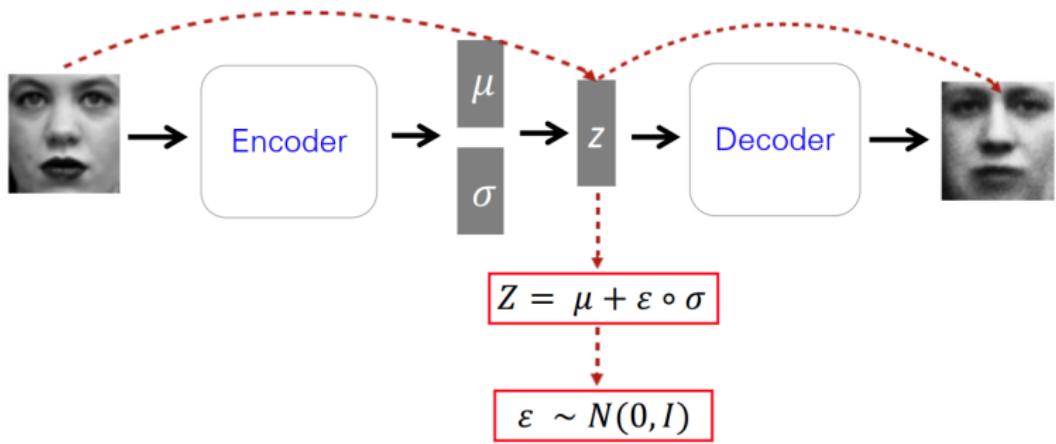
# Reparametrization Trick



# Reparametrization Trick



# Reparametrization Trick



# Training VAE

**Traditional AE:**

Input Image:



Output Images:



**Variational AE:**

Input Image:



Output Images:



Difference:



# More about Generative models

Next Class